

## APPENDIX

### A. Proof of Theorem 1

Let us rewrite the risks and generalization error under the DNN setup. Let  $(X, Y) \sim P_{X,Y}$  be a pair of test data sample. At each layer  $l$ , the internal representation  $T_l$  of a test data feature  $X$  is conditionally independent of  $W_{l+1}^L$  given  $W_1^l$ . For any  $\mathbf{W} \in \mathcal{W}$ , let the loss function be rewritten as  $\ell(\mathbf{W}, X, Y) = \ell(g_{\mathbf{W}_L} \circ g_{\mathbf{W}_{L-1}} \circ \cdots \circ g_{\mathbf{W}_1}(X), Y)$ . The expected population risk over all possible  $\mathbf{W}$  is given by

$$\mathbb{E}_{\mathbf{W}}[\mathcal{L}_P(\mathbf{W}, P_{X,Y})] = \mathbb{E}[\mathbb{E}[\ell(g_{\mathbf{W}_L} \circ g_{\mathbf{W}_{L-1}} \circ \cdots \circ g_{\mathbf{W}_{l+1}}(T_l), Y) | \mathbf{W}_1^l]]$$

where  $l \in [L]$  and given  $\mathbf{W}_1^l$ ,  $(T_l, Y)$  are independent of  $\mathbf{W}_{l+1}^L$ .

Denote the overall feature mapping function as  $f_{\mathbf{W}} \triangleq g_{\mathbf{W}_L} \circ g_{\mathbf{W}_{L-1}} \circ \cdots \circ g_{\mathbf{W}_1}$ . Similarly, for any  $l \in [L]$ , the expected empirical risk can also be rewritten as

$$\mathbb{E}[\mathcal{L}_E(\mathbf{W}, D_n)] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbb{E}[\ell(g_{\mathbf{W}_L} \circ g_{\mathbf{W}_{L-1}} \circ \cdots \circ g_{\mathbf{W}_{l+1}}(T_{l,i}), Y_i) | \mathbf{W}_1^l]].$$

For notational simplicity, let  $g_{\mathbf{W}_k^j} := g_{\mathbf{W}_k} \circ g_{\mathbf{W}_{k-1}} \circ \cdots \circ g_{\mathbf{W}_j}$  for any  $k < j$  and  $k, j \in \mathbb{N}$ . Then the expected generalization error can be rewritten as

$$\text{gen}(P_{\mathbf{W}|D_n}, P_{X,Y}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[ \mathbb{E}[\ell(g_{\mathbf{W}_{l+1}^L}(T_l), Y) | \mathbf{W}_1^l] - \mathbb{E}[\ell(g_{\mathbf{W}_{l+1}^L}(T_{l,i}), Y_i) | \mathbf{W}_1^l] \right]. \quad (2)$$

If the loss function  $\ell(\mathbf{w}, X, Y)$  is  $\sigma$ -sub-Gaussian under  $P_{X,Y}$  for all  $\mathbf{w} \in \mathcal{W}$ , we also have for any  $l \in [0 : L]$ ,  $\ell(g_{\mathbf{W}_{l+1}^L}(T_l), Y)$  is  $\sigma$ -sub-Gaussian under  $P_{T_l,Y|\mathbf{W}=\mathbf{w}}$  for all  $\mathbf{w} \in \mathcal{W}$ . From Donsker-Varadhan representation, we have for any  $\lambda \in \mathbb{R}$ ,

$$\begin{aligned} & D_{\text{KL}}(P_{\mathbf{W}_{l+1}^L, T_{l,i}, Y_i | \mathbf{W}_1^l} \| P_{T_l, Y | \mathbf{W}_1^l} \otimes P_{\mathbf{W}_{l+1}^L | \mathbf{W}_1^l}) \\ & \geq \mathbb{E}_{\mathbf{W}_{l+1}^L, T_{l,i}, Y_i | \mathbf{W}_1^l} [\lambda \ell(g_{\mathbf{W}_{l+1}^L}(T_{l,i}), Y_i)] - \log \mathbb{E}_{\mathbf{W}_{l+1}^L | \mathbf{W}_1^l} \mathbb{E}_{T_l, Y | \mathbf{W}_1^l} [\exp(\lambda \ell(g_{\mathbf{W}_{l+1}^L}(T_l), Y))] \\ & \geq \lambda (\mathbb{E}_{\mathbf{W}_{l+1}^L, T_{l,i}, Y_i | \mathbf{W}_1^l} [\ell(g_{\mathbf{W}_{l+1}^L}(T_{l,i}), Y_i)] - \mathbb{E}_{\mathbf{W}_{l+1}^L | \mathbf{W}_1^l} \mathbb{E}_{T_l, Y | \mathbf{W}_1^l} [\ell(g_{\mathbf{W}_{l+1}^L}(T_l), Y)]) - \frac{\lambda^2 \sigma^2}{2}. \end{aligned}$$

We can decompose  $D_{\text{KL}}(P_{\mathbf{W}_{l+1}^L, T_{l,i}, Y_i | \mathbf{W}_1^l} \| P_{T_l, Y | \mathbf{W}_1^l} \otimes P_{\mathbf{W}_{l+1}^L | \mathbf{W}_1^l} | P_{\mathbf{W}_1^l})$  as follows

$$\begin{aligned} & D_{\text{KL}}(P_{\mathbf{W}_{l+1}^L, T_{l,i}, Y_i | \mathbf{W}_1^l} \| P_{T_l, Y | \mathbf{W}_1^l} \otimes P_{\mathbf{W}_{l+1}^L | \mathbf{W}_1^l} | P_{\mathbf{W}_1^l}) \\ & = D_{\text{KL}}(P_{\mathbf{W}_{l+1}^L, T_{l,i}, Y_i | \mathbf{W}_1^l} \| P_{T_{l,i}, Y_i | \mathbf{W}_1^l} \otimes P_{\mathbf{W}_{l+1}^L | \mathbf{W}_1^l} | P_{\mathbf{W}_1^l}) + D_{\text{KL}}(P_{T_{l,i}, Y_i | \mathbf{W}_1^l} \| P_{T_l, Y | \mathbf{W}_1^l} | P_{\mathbf{W}_1^l}) \\ & = I(T_{l,i}, Y_i; \mathbf{W}_{l+1}^L | \mathbf{W}_1^l) + D_{\text{KL}}(P_{T_{l,i}, Y_i | \mathbf{W}_1^l} \| P_{T_l, Y | \mathbf{W}_1^l} | P_{\mathbf{W}_1^l}). \end{aligned} \quad (3)$$

Thus, we have

$$\begin{aligned} & I(T_{l,i}, Y_i; \mathbf{W}_{l+1}^L | \mathbf{W}_1^l) + D_{\text{KL}}(P_{T_{l,i}, Y_i | \mathbf{W}_1^l} \| P_{T_l, Y | \mathbf{W}_1^l} | P_{\mathbf{W}_1^l}) = D_{\text{KL}}(P_{\mathbf{W}_{l+1}^L, T_{l,i}, Y_i | \mathbf{W}_1^l} \| P_{T_l, Y | \mathbf{W}_1^l} \otimes P_{\mathbf{W}_{l+1}^L | \mathbf{W}_1^l} | P_{\mathbf{W}_1^l}) \\ & \geq \lambda \mathbb{E}_{\mathbf{W}_1^l} [\mathbb{E}_{\mathbf{W}_{l+1}^L, T_{l,i}, Y_i | \mathbf{W}_1^l} [\ell(g_{\mathbf{W}_{l+1}^L}(T_{l,i}), Y_i)] - \mathbb{E}_{\mathbf{W}_{l+1}^L | \mathbf{W}_1^l} \mathbb{E}_{T_l, Y | \mathbf{W}_1^l} [\ell(g_{\mathbf{W}_{l+1}^L}(T_l), Y)]] - \frac{\lambda^2 \sigma^2}{2}. \end{aligned}$$

By optimizing the RHS over  $\lambda > 0$  and  $\lambda \leq 0$ , respectively, we finally obtain

$$\begin{aligned} & \left| \mathbb{E}_{\mathbf{W}_1^l} \mathbb{E}_{\mathbf{W}_{l+1}^L, T_{l,i}, Y_i | \mathbf{W}_1^l} [\ell(g_{\mathbf{W}_{l+1}^L}(T_{l,i}), Y_i)] - \mathbb{E}_{\mathbf{W}_1^l} \mathbb{E}_{\mathbf{W}_{l+1}^L | \mathbf{W}_1^l} \mathbb{E}_{T_l, Y | \mathbf{W}_1^l} [\ell(g_{\mathbf{W}_{l+1}^L}(T_l), Y)] \right| \\ & \leq \sqrt{2\sigma^2 (I(T_{l,i}, Y_i; \mathbf{W}_{l+1}^L | \mathbf{W}_1^l) + D_{\text{KL}}(P_{T_{l,i}, Y_i | \mathbf{W}_1^l} \| P_{T_l, Y | \mathbf{W}_1^l} | P_{\mathbf{W}_1^l}))}, \end{aligned}$$

which holds for all  $l \in [L]$ . Conditioned on  $\mathbf{W}_l$ ,  $T_{l,i}$  and  $T_l$  are generated by the same process from  $T_{l-1,i}$  and  $T_{l-1}$ , respectively. By the data-processing inequality, the KL divergence in (3) can be bounded as follows:

$$\begin{aligned} D_{\text{KL}}(P_{\mathbf{W}_{l+1}^L, T_{l,i}, Y_i | \mathbf{W}_1^l} \| P_{T_l, Y | \mathbf{W}_1^l} \otimes P_{\mathbf{W}_{l+1}^L | \mathbf{W}_1^l} | P_{\mathbf{W}_1^l}) & \leq D_{\text{KL}}(P_{\mathbf{W}_{l+1}^L, T_{l-1,i}, Y_i | \mathbf{W}_1^l} \| P_{T_{l-1}, Y | \mathbf{W}_1^l} \otimes P_{\mathbf{W}_{l+1}^L | \mathbf{W}_1^l} | P_{\mathbf{W}_1^l}) \\ & = D_{\text{KL}}(P_{\mathbf{W}_l^L, T_{l-1,i}, Y_i | \mathbf{W}_1^{l-1}} \| P_{T_{l-1}, Y | \mathbf{W}_1^{l-1}} \otimes P_{\mathbf{W}_l^L | \mathbf{W}_1^{l-1}} | P_{\mathbf{W}_1^{l-1}}) \\ & \vdots \\ & \leq D_{\text{KL}}(P_{X_i, Y_i, \mathbf{W}_1^L} \| P_{X,Y} \otimes P_{\mathbf{W}_1^L}) = I(X_i, Y_i; \mathbf{W}). \end{aligned}$$

Theorem 1 can be thus proved by induction.

### B. Proof of Theorem 2

Recall the Kantorovich-Rubinstein duality [36]: for any two probability measures  $P, Q \in \mathcal{P}(\mathcal{X})$ ,  $W_1(P, Q) = \sup_{f \in \text{Lip}_1(\mathcal{X})} \mathbb{E}_P[f] - \mathbb{E}_Q[f]$ , where  $\text{Lip}_k(\mathcal{X}) = \{f \in \{f : \mathcal{X} \rightarrow \mathbb{R}\} : |f(x) - f(y)| \leq k\|x - y\|, \forall x, y \in \mathcal{X}\}$ , for any  $k \in \mathbb{R}_{\geq 0}$ .

Since  $\tilde{\ell}(g_{\mathbf{w}_L} \circ \dots \circ g_{\mathbf{w}_1}(X), Y)$  is  $\rho_0$ -Lipschitz in  $(g_{\mathbf{w}_L} \circ \dots \circ g_{\mathbf{w}_1}(X), Y)$  and  $\phi_l(\cdot)$  is  $\rho_l$ -Lipschitz, we have for any  $\mathbf{w}$ ,

$$\begin{aligned} |\ell(\mathbf{w}, x, y) - \ell(\mathbf{w}, x', y')| &= |\tilde{\ell}(g_{\mathbf{w}_L} \circ \dots \circ g_{\mathbf{w}_1}(x), y) - \tilde{\ell}(g_{\mathbf{w}_L} \circ \dots \circ g_{\mathbf{w}_1}(x'), y')| \\ &\leq \rho_0 \|(g_{\mathbf{w}_L}(g_{\mathbf{w}_1^{L-1}}(x)), y) - (g_{\mathbf{w}_L}(g_{\mathbf{w}_1^{L-1}}(x')), y')\| \\ &\leq \rho_0 \sqrt{(\rho_L \|\mathbf{w}_L\| \|g_{\mathbf{w}_1^{L-1}}(x) - g_{\mathbf{w}_1^{L-1}}(x')\|)^2 + (y - y')^2} \\ &\vdots \\ &\leq \bar{\rho}_0(\mathbf{w}) \sqrt{\|x - x'\|^2 + (y - y')^2} \end{aligned}$$

where  $\bar{\rho}_0(\mathbf{w}) := \rho_0(1 \vee \prod_{j=1}^L \rho_j \|\mathbf{w}_j\|)$ . It means  $\ell(\mathbf{W}, X, Y)$  is  $\bar{\rho}_0(\mathbf{W})$ -Lipschitz in  $(X, Y)$  for any  $\mathbf{W}$ . Then we have

$$\text{gen}(P_{\mathbf{W}|D_n}, P_{X,Y}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\ell(\mathbf{W}, X, Y) - \ell(\mathbf{W}, X_i, Y_i)] \leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\bar{\rho}_0(\mathbf{W}) W_1(P_{X_i, Y_i | \mathbf{W}}, P_{X, Y | \mathbf{W}})].$$

For  $l = 1, \dots, L$ , similarly, we have  $\tilde{\ell}(g_{\mathbf{w}_L} \circ \dots \circ g_{\mathbf{w}_{l+1}}(T_l), Y)$  is  $\rho_0(1 \vee \prod_{j=l+1}^L \rho_j \|\mathbf{w}_j\|)$ -Lipschitz in  $(T_l, Y)$  for all  $\mathbf{W}$ . Let  $\bar{\rho}_l(\mathbf{W}) = \rho_0(1 \vee \prod_{j=l+1}^L \rho_j \|\mathbf{w}_j\|)$ . Then from the definition in (2), we have

$$\text{gen}(P_{\mathbf{W}|D_n}, P_{X,Y}) \leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\bar{\rho}_l(\mathbf{W}) W_1(P_{T_{l,i}, Y_i | \mathbf{W}}, P_{T_l, Y | \mathbf{W}})].$$

The proof is completed by taking the minimum over  $l = 0, \dots, L$ .

### C. Proof of Remark 1

Let  $\text{diam}(\mathcal{X}) := \sup\{\|x - y\| : x, y \in \mathcal{X}\}$ . From [35, Theorem 4], Pinsker's and Bretagnolle-Huber inequalities, for any two probability distributions  $\mu, \nu \in \mathcal{P}(\mathcal{X})$ , we have

$$W_1(\mu, \nu) \leq \text{diam}(\mathcal{X}) D_{\text{TV}}(\mu, \nu) \leq \text{diam}(\mathcal{X}) \sqrt{\left(\frac{1}{2} D_{\text{KL}}(\mu \| \nu) \wedge (1 - \exp(-D_{\text{KL}}(\mu \| \nu)))\right)}.$$

From Theorem 1 and [37], the generalization error can be bounded as follows:

$$|\text{gen}(P_{\mathbf{W}|D_n}, P_{X,Y})| \leq \frac{A}{n} \sum_{i=1}^n D_{\text{TV}}(P_{T_{L,i}, Y_i | \mathbf{W}}, P_{T_L, Y | \mathbf{W}} | P_{\mathbf{W}}) \leq \text{UB}(L),$$

where  $\text{UB}(L) := \frac{\sqrt{2}A}{2n} \sum_{i=1}^n \sqrt{D_{\text{KL}}(P_{T_{L,i}, Y_i | \mathbf{W}} \| P_{T_L, Y | \mathbf{W}} | P_{\mathbf{W}})}$ . It can be observed that the total variation distance based bound is tighter under this condition.

Let  $l^* := \min_{l=0, \dots, L} \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\bar{\rho}_l(\mathbf{W}) W_1(P_{T_{l,i}, Y_i | \mathbf{W}}, P_{T_l, Y | \mathbf{W}})]$ .

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\bar{\rho}_{l^*}(\mathbf{W}) W_1(P_{T_{l^*,i}, Y_i | \mathbf{W}}, P_{T_{l^*}, Y | \mathbf{W}})] &\leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\bar{\rho}_L(\mathbf{W}) W_1(P_{T_{L,i}, Y_i | \mathbf{W}}, P_{T_L, Y | \mathbf{W}})] \\ &= \frac{\rho_0}{n} \sum_{i=1}^n W_1(P_{T_{L,i}, Y_i | \mathbf{W}}, P_{T_L, Y | \mathbf{W}} | P_{\mathbf{W}}) \leq \frac{\rho_0 \text{diam}(\mathcal{T}_L \times \mathcal{Y})}{n} \sum_{i=1}^n D_{\text{TV}}(P_{T_{L,i}, Y_i | \mathbf{W}}, P_{T_L, Y | \mathbf{W}} | P_{\mathbf{W}}). \end{aligned}$$

Under our supervised classification setting,  $\text{diam}(\mathcal{T}_L \times \mathcal{Y}) = K^2$ . Thus, when  $\rho_0 K^2 \leq A$ , 1-Wasserstein distance based bound is even tighter than the one based on the TV-distance.

#### D. Proofs for the case study of binary Gaussian mixture classification

To simplify some of the notation ahead, the distribution of a Gaussian random variable  $X$  with mean  $\mu$  and variance  $\sigma^2$  is denoted by  $\mathcal{N}_X(\mu, \sigma^2)$ . Under the binary Gaussian mixture classification setting, we first know that the prior of  $\mathbf{W}_{\otimes L}^\top$  is  $P_{\mathbf{W}_{\otimes L}^\top} = \mathcal{N}_{\mathbf{W}_{\otimes L}}(\mu_0, \frac{\sigma_0^2}{n} \mathbf{I}_{d_0})$ . Given any pair of training data sample  $(X_i, Y_i)$ , we have

$$\mathbf{W}_{\otimes L}^\top | (X_i, Y_i) = \frac{1}{n} Y_i X_i + \frac{1}{n} \sum_{j \neq i}^n Y_j X_j \sim \mathcal{N}_{\mathbf{W}_{\otimes L}}(\mu_{\mathbf{W}_{\otimes L}|i}, \Sigma_{\mathbf{W}_{\otimes L}|i}),$$

where  $\mu_{\mathbf{W}_{\otimes L}|i} = \frac{1}{n} Y_i X_i + \frac{n-1}{n} \mu_0$  and  $\Sigma_{\mathbf{W}_{\otimes L}|i} = \frac{n-1}{n^2} \sigma_0^2 \mathbf{I}_{d_0}$ . Then the posterior distribution of  $(X_i, Y_i)$  given  $\mathbf{W}_{\otimes L}$  is given by

$$\begin{aligned} P_{X_i, Y_i | \mathbf{W}_{\otimes L}} &= \frac{P_{\mathbf{W}_{\otimes L} | X_i, Y_i} P_{X_i, Y_i}}{P_{\mathbf{W}_{\otimes L}}} = \frac{\mathcal{N}_{\mathbf{W}_{\otimes L}}(\mu_{\mathbf{W}_{\otimes L}|i}, \Sigma_{\mathbf{W}_{\otimes L}|i})}{\mathcal{N}_{\mathbf{W}_{\otimes L}}(\mu_0, \frac{\sigma_0^2}{n} \mathbf{I}_{d_0})} \times \frac{1}{2} \mathcal{N}_{X_i}(Y_i \mu_0, \sigma_0^2 \mathbf{I}_{d_0}) \\ &= \frac{1}{2} \mathcal{N}_{X_i}(Y_i \mu_0, \sigma_0^2 \mathbf{I}_{d_0}) \times C_i \mathcal{N}_{\mathbf{W}_{\otimes L}}\left(Y_i X_i, \frac{(n-1)\sigma_0^2}{n} \mathbf{I}_{d_0}\right) = \frac{1}{2} \mathcal{N}_{Y_i X_i}\left(\mathbf{W}_{\otimes L}^\top, \frac{(n-1)\sigma_0^2}{n} \mathbf{I}_{d_0}\right), \end{aligned}$$

where  $C_i = n^{d_0} \sqrt{\left(\frac{2\pi\sigma_0^2}{n^2}\right)^{d_0}} \exp\left\{\frac{1}{2\sigma_0^2} (Y_i X_i - \mu_0)^\top (Y_i X_i - \mu_0)\right\}$ . By integrating  $P_{X_i, Y_i | \mathbf{W}_{\otimes L}}$  over  $X_i$ , we obtain  $P_{Y_i | \mathbf{W}_{\otimes L}} = \frac{1}{2}$ . We can also conclude that  $P_{X_i, Y_i | \mathbf{W}, \mathbf{W}_{\otimes L}} = P_{X_i, Y_i | \mathbf{W}_{\otimes L}}$  and  $P_{Y_i | \mathbf{W}, \mathbf{W}_{\otimes L}} = P_{Y_i | \mathbf{W}_{\otimes L}} = \frac{1}{2}$ .

Next, we compute the divergences between the prior and posterior.

*Proof of Proposition 4.* For any  $l \in [L]$ , conditioned on  $(Y_i, \mathbf{W}, \mathbf{W}_{\otimes L})$ , the distribution of  $T_{l,i} = \mathbf{W}_{\otimes l} X_i$  is Gaussian with the mean and covariance

$$\mathbb{E}[T_{l,i} | Y_i, \mathbf{W}, \mathbf{W}_{\otimes L}] = Y_i \mathbf{W}_{\otimes l} \mathbf{W}_{\otimes L}^\top, \quad \text{Cov}[T_{l,i} | Y_i, \mathbf{W}, \mathbf{W}_{\otimes L}] = \frac{(n-1)\sigma_0^2}{n} \mathbf{W}_{\otimes l} \mathbf{W}_{\otimes L}^\top.$$

Similarly for a test data sample, when conditioned on  $(Y, \mathbf{W}, \mathbf{W}_{\otimes L})$ , the distribution of  $T_l = \mathbf{W}_{\otimes l} X$  is Gaussian with mean and covariance

$$\mathbb{E}[T_l | Y, \mathbf{W}, \mathbf{W}_{\otimes L}] = \mathbb{E}[T_l | Y, \mathbf{W}] = Y \mathbf{W}_{\otimes l} \mu_0, \quad \text{Cov}[T_l | Y, \mathbf{W}, \mathbf{W}_{\otimes L}] = \sigma_0^2 \mathbf{W}_{\otimes l} \mathbf{W}_{\otimes L}^\top.$$

Note that when  $\mathbf{W}_{\otimes l}$  is not a full-rank matrix, the covariance matrices  $\text{Cov}[T_{l,i} | Y_i, \mathbf{W}, \mathbf{W}_{\otimes L}]$  and  $\text{Cov}[T_l | Y, \mathbf{W}, \mathbf{W}_{\otimes L}]$  are singular. Thus, the posteriors  $P_{T_{l,i} | Y_i, \mathbf{W}, \mathbf{W}_{\otimes L}}$  and  $P_{T_l | Y, \mathbf{W}, \mathbf{W}_{\otimes L}}$  are the push-forwards of two  $d_l$ -dimensional nonsingular Gaussian distributions to the lower-dimensional space. In fact, the dimension can be further proven to be  $\text{rank}(\mathbf{W}_{\otimes l})$  via eigendecomposition. Take the test data sample  $(X, Y = 1)$  for example in the followings. Let  $X = \mu_0 + Z$  and the  $l^{\text{th}}$  representation  $T_l = \mathbf{W}_{\otimes l} X = \mathbf{W}_{\otimes l} \mu_0 + \mathbf{W}_{\otimes l} Z$ , where  $Z \sim \mathcal{N}(0, \sigma_0^2 \mathbf{I}_{d_0})$  and  $\mathbf{W}_{\otimes l} Z \sim \mathcal{N}(0, \sigma_0^2 \mathbf{W}_{\otimes l} \mathbf{W}_{\otimes L}^\top)$  is a singular Gaussian distribution. Since  $\mathbf{W}_{\otimes l} \mathbf{W}_{\otimes L}^\top$  is a  $(d_l \times d_l)$  positive semi-definite and symmetric matrix, the eigendecomposition is given by  $\mathbf{W}_{\otimes l} \mathbf{W}_{\otimes L}^\top = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$ , where  $\mathbf{U}$  is a  $(d_l \times d_l)$  orthogonal matrix,  $\mathbf{\Lambda}$  is a  $(d_l \times d_l)$  diagonal matrix whose entries are the eigenvalues of  $\mathbf{W}_{\otimes l} \mathbf{W}_{\otimes L}^\top$ . We assume that  $\mathbf{W}_{\otimes l}$  satisfies that only the first  $\text{rank}(\mathbf{W}_{\otimes l})$  entries of  $\mathbf{\Lambda}$  are non-zero. Construct the following two  $d_l$ -dimensional vectors:

$$Z_0 := (Z, \underbrace{0, \dots, 0}_{(d_l - d_0) \text{ entries}})^\top, \quad \tilde{Z} := (\underbrace{\tilde{Z}_1, \dots, \tilde{Z}_{\text{rank}(\mathbf{W}_{\otimes l})}}_{\sim \mathcal{N}(0, \sigma_0^2 \mathbf{I}_{\text{rank}(\mathbf{W}_{\otimes l})})}, \underbrace{0, \dots, 0}_{(d_l - \text{rank}(\mathbf{W}_{\otimes l})) \text{ entries}})^\top.$$

Since the last  $(d_l - \text{rank}(\mathbf{W}_{\otimes l}))$  columns of  $\mathbf{U} \mathbf{\Lambda}^{\frac{1}{2}}$  are all zero, the following equalities hold:

$$\mathbf{W}_{\otimes l} Z \stackrel{d}{=} \mathbf{U} \mathbf{\Lambda}^{\frac{1}{2}} Z_0 \stackrel{d}{=} \mathbf{U} \mathbf{\Lambda}^{\frac{1}{2}} \tilde{Z}, \quad \text{and} \quad T_l \stackrel{d}{=} \mathbf{W}_{\otimes l} \mu_0 + \mathbf{U} \mathbf{\Lambda}^{\frac{1}{2}} \tilde{Z}.$$

Thus, only the a subset of  $\text{rank}(\mathbf{W}_{\otimes l})$  covariates of  $T_l$  are effective random variables and the Gaussian PDF of  $T_l$  is on the  $\text{rank}(\mathbf{W}_{\otimes l})$ -dimensional space. Let  $r_l = \text{rank}(\mathbf{W}_{\otimes l})$ . We can define a restriction of Lebesgue measure to the  $\text{rank}(\mathbf{W}_{\otimes l})$ -dimensional affine subspace of  $\mathbb{R}^{r_l}$  where the Gaussian distribution is supported. With respect to this measure the distribution of  $T_l$  given  $(\mathbf{W}, \mathbf{W}_{\otimes l}, Y = 1)$  has the density of the following motif:

$$p_{\mathbf{W}, \mathbf{W}_{\otimes l}, Y=1}(T_l) = \frac{\exp\left(-\frac{1}{2\sigma_0^2} (T_l - \mathbf{W}_{\otimes l} \mu_0)^\top (\mathbf{W}_{\otimes l} \mathbf{W}_{\otimes L}^\top)^\dagger (T_l - \mathbf{W}_{\otimes l} \mu_0)\right)}{\sqrt{(2\pi)^{r_l} \det^*(\sigma_0^2 \mathbf{W}_{\otimes l} \mathbf{W}_{\otimes L}^\top)}} \quad (4)$$

where  $(\mathbf{W}_{\otimes l} \mathbf{W}_{\otimes L}^\top)^\dagger = \mathbf{U} \mathbf{\Lambda}^\dagger \mathbf{U}^\top$  is the generalized inverse of  $\mathbf{W}_{\otimes l} \mathbf{W}_{\otimes L}^\top$ ,  $\mathbf{\Lambda}^\dagger$  is the pseudo-inverse of  $\mathbf{\Lambda}$ , and  $\det^*$  is the pseudo-determinant. In a similar manner, the density of  $T_{l,i}$  can be obtained by replacing the mean  $\mathbf{W}_{\otimes l} \mu_0$  with  $\mathbf{W}_{\otimes l} \mathbf{W}_{\otimes L}^\top$  and  $\sigma_0^2$  with  $\frac{(n-1)\sigma_0^2}{n}$ .

Recall that the KL divergence between any two Gaussian distributions  $P = \mathcal{N}(\mu_p, \Sigma_p), Q = \mathcal{N}(\mu_q, \Sigma_q) \in \mathcal{P}(\mathbb{R}^k)$  for some  $k \in \mathbb{N}$  is given by

$$D_{\text{KL}}(P\|Q) = \frac{1}{2} \left( \log \frac{\det^*(\Sigma_q)}{\det^*(\Sigma_p)} - k + (\mu_p - \mu_q)^\top \Sigma_q^\dagger (\mu_p - \mu_q) + \text{tr}(\Sigma_q^\dagger \Sigma_p) \right).$$

For  $l = 1, \dots, L$ , the  $\text{UB}(l)$  in Theorem 1 can be written as

$$\text{UB}(l) = \frac{\sqrt{2}\sigma}{n} \sum_{i=1}^n \sqrt{D_{\text{KL}}(P_{T_{l,i}, Y_i} | \mathbf{W} \| P_{T_l, Y} | \mathbf{W} | P_{\mathbf{W}})} \leq \frac{\sqrt{2}\sigma}{n} \sum_{i=1}^n \sqrt{D_{\text{KL}}(P_{T_{l,i}, Y_i} | \mathbf{W}, \mathbf{W}_{\otimes L} \| P_{T_l, Y} | \mathbf{W}, \mathbf{W}_{\otimes L} | P_{\mathbf{W}, \mathbf{W}_{\otimes L}})} =: \widetilde{\text{UB}}(l),$$

where the hierarchical structure  $\widetilde{\text{UB}}(L) \leq \widetilde{\text{UB}}(L-1) \leq \dots \leq \widetilde{\text{UB}}(0)$  still holds.

From the probability density function (PDF) of  $T_l$  and  $T_{l,i}$  (c.f. (4)), the KL divergence term in the upper bound  $\widetilde{\text{UB}}(l)$  can be rewritten as: for  $l = 1, \dots, L$ ,

$$\begin{aligned} & D_{\text{KL}}(P_{T_{l,i}, Y_i} | \mathbf{W}, \mathbf{W}_{\otimes L} \| P_{T_l, Y} | \mathbf{W}, \mathbf{W}_{\otimes L} | P_{\mathbf{W}, \mathbf{W}_{\otimes L}}) \\ &= \frac{1}{2} D_{\text{KL}}(P_{T_{l,i} | Y_i=1} | \mathbf{W}, \mathbf{W}_{\otimes L} \| P_{T_l | Y=1} | \mathbf{W}, \mathbf{W}_{\otimes L} | P_{\mathbf{W}, \mathbf{W}_{\otimes L}}) + \frac{1}{2} D_{\text{KL}}(P_{T_{l,i} | Y_i=-1} | \mathbf{W}, \mathbf{W}_{\otimes L} \| P_{T_l | Y=-1} | \mathbf{W}, \mathbf{W}_{\otimes L} | P_{\mathbf{W}, \mathbf{W}_{\otimes L}}) \\ &= \frac{1}{2} \mathbb{E} \left[ \mathbb{E} \left[ r_l \left( \log \frac{n}{n-1} - 1 + \frac{n-1}{n} \right) + \frac{1}{\sigma_0^2} (\mu_0 - \mathbf{W}_{\otimes L}^\top)^\top (\mathbf{W}_{\otimes L})^\top (\mathbf{W}_{\otimes L} \mathbf{W}_{\otimes L}^\top)^\dagger \mathbf{W}_{\otimes L} (\mu_0 - \mathbf{W}_{\otimes L}^\top) \middle| \mathbf{W}, \mathbf{W}_{\otimes L} \right] \right] \\ &= \frac{1}{2} \mathbb{E} \left[ \mathbb{E} \left[ r_l \left( \log \frac{n}{n-1} - \frac{1}{n} \right) + \frac{1}{\sigma_0^2} (\mu_0 - \mathbf{W}_{\otimes L}^\top)^\top \mathbf{V} \mathbf{S}^\top \mathbf{U}^\top \mathbf{U} \Lambda^\dagger \mathbf{U}^\top \mathbf{U} \mathbf{S} \mathbf{V}^\top (\mu_0 - \mathbf{W}_{\otimes L}^\top) \middle| \mathbf{W}, \mathbf{W}_{\otimes L} \right] \right] \\ &= \mathbb{E} \left[ \frac{r_l}{2} \left( \log \frac{n}{n-1} - \frac{1}{n} \right) + \frac{1}{2\sigma_0^2} \mathbb{E}[\|\mu_0 - \mathbf{W}_{\otimes L}^\top\|^2] \right] \\ &= \frac{\mathbb{E}[r_l]}{2} \left( \log \frac{n}{n-1} - \frac{1}{n} \right) + \frac{d_0}{2n}, \end{aligned}$$

where the last equality follows since  $\frac{\sqrt{n}}{\sigma_0}(\mu_0 - \mathbf{W}_{\otimes L}^\top) \sim \mathcal{N}(0, \mathbf{I}_{d_0})$ ,  $\frac{n}{\sigma_0^2} \|\mu_0 - \mathbf{W}_{\otimes L}^\top\|^2 \sim \chi_{d_0}^2$  and  $\frac{1}{2\sigma_0^2} \mathbb{E}[\|\mu_0 - \mathbf{W}_{\otimes L}^\top\|^2] = \frac{d_0}{2n}$ .

The KL divergence term in the upper bound  $\widetilde{\text{UB}}(0)$  can be rewritten as

$$\begin{aligned} D_{\text{KL}}(P_{X_i, Y_i} | \mathbf{W}, \mathbf{W}_{\otimes L} \| P_{X, Y} | \mathbf{W}, \mathbf{W}_{\otimes L}) &= \mathbb{E} \left[ \mathbb{E} \left[ D_{\text{KL}} \left( \frac{1}{2} \mathcal{N}_{X_i} \left( -\mathbf{W}_{\otimes L}^\top, \frac{(n-1)\sigma_0^2}{n} \mathbf{I}_{d_0} \right) \middle\| \frac{1}{2} \mathcal{N}_X(-\mu_0, \sigma_0^2 \mathbf{I}_{d_0}) \right) \right. \right. \\ &\quad \left. \left. + D_{\text{KL}} \left( \frac{1}{2} \mathcal{N}_{X_i} \left( \mathbf{W}_{\otimes L}^\top, \frac{(n-1)\sigma_0^2}{n} \mathbf{I}_{d_0} \right) \middle\| \frac{1}{2} \mathcal{N}_X(\mu_0, \sigma_0^2 \mathbf{I}_{d_0}) \right) \middle| \mathbf{W}, \mathbf{W}_{\otimes L} \right] \right] \\ &= \frac{1}{2} \mathbb{E} \left[ \mathbb{E} \left[ d_0 \left( \log \frac{n}{n-1} - 1 + \frac{n-1}{n} \right) + \frac{1}{\sigma_0^2} \|\mu_0 - \mathbf{W}_{\otimes L}^\top\|^2 \middle| \mathbf{W}, \mathbf{W}_{\otimes L} \right] \right] \\ &= \frac{d_0}{2} \left( \log \frac{n}{n-1} - \frac{1}{n} \right) + \frac{1}{2\sigma_0^2} \mathbb{E}[\|\mu_0 - \mathbf{W}_{\otimes L}^\top\|^2] \\ &= \frac{d_0}{2} \left( \log \frac{n}{n-1} - \frac{1}{n} \right) + \frac{d_0}{2n}. \end{aligned}$$

□

*Proof of Proposition 5.* Since the closed form of  $W_1$  between two Gaussian distributions is not known but known for  $W_2$  and  $W_1(\cdot, \cdot) \leq W_2(\cdot, \cdot)$ , we consider analysing  $W_2(P_{T_{l,i}, Y_i} | \mathbf{W}, P_{T_l, Y} | \mathbf{W} | P_{\mathbf{W}})$  as a surrogate of the upper bound in Theorem 2. Following the proof of Theorem 2, we can obtain

$$\begin{aligned} \text{gen}(P_{\mathbf{W} | D_n}, P_{X, Y}) &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\ell(\mathbf{W}, X, Y) - \ell(\mathbf{W}, X_i, Y_i)] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbb{E}[\ell(\mathbf{W}, X, Y) - \ell(\mathbf{W}, X_i, Y_i)] | \mathbf{W}, \mathbf{W}_{\otimes L}, Y_i] \\ &\stackrel{(a)}{\leq} \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\bar{\rho}_0(\mathbf{W}) W_1(P_{X_i | Y_i, \mathbf{W}, \mathbf{W}_{\otimes L}}, P_{X | Y, \mathbf{W}, \mathbf{W}_{\otimes L}})] \\ &\leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\bar{\rho}_0(\mathbf{W}) W_2(P_{X_i | Y_i, \mathbf{W}, \mathbf{W}_{\otimes L}}, P_{X | Y, \mathbf{W}, \mathbf{W}_{\otimes L}})], \end{aligned}$$

where (a) follows since  $P_{Y_i|\mathbf{W}, \mathbf{W}_{\otimes L}} = P_Y = \text{Unif}\{-1, +1\}$ . Similarly, we also have for all  $l = 1, \dots, L$ .

$$\text{gen}(P_{\mathbf{W}|D_n}, P_{X,Y}) \leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\bar{\rho}_l(\mathbf{W}) \mathcal{W}_2(P_{T_{l,i}|Y_i, \mathbf{W}, \mathbf{W}_{\otimes L}}, P_{T_l|Y, \mathbf{W}, \mathbf{W}_{\otimes L}})].$$

By plugging the  $P_{T_{l,i}|Y_i, \mathbf{W}, \mathbf{W}_{\otimes L}}$  and  $P_{T_l|Y, \mathbf{W}, \mathbf{W}_{\otimes L}}$  (c.f. (4)) into the upper bound, we have

$$\begin{aligned} \mathcal{W}_2(P_{T_{l,i}|Y_i, \mathbf{W}, \mathbf{W}_{\otimes L}}, P_{T_l|Y, \mathbf{W}, \mathbf{W}_{\otimes L}}) &= \left( \|\mathbf{W}_{\otimes l}(\mathbf{W}_{\otimes L}^\top - \mu_0)\|^2 + \text{tr} \left( \frac{(n-1)\sigma_0^2}{n} (\mathbf{W}_{\otimes l} \mathbf{W}_{\otimes l}^\top) + \sigma_0^2 (\mathbf{W}_{\otimes l} \mathbf{W}_{\otimes l}^\top) \right. \right. \\ &\quad \left. \left. - 2 \left( \frac{(n-1)\sigma_0^4}{n} (\mathbf{W}_{\otimes l} \mathbf{W}_{\otimes l}^\top)^{\frac{1}{2}} (\mathbf{W}_{\otimes l} \mathbf{W}_{\otimes l}^\top) (\mathbf{W}_{\otimes l} \mathbf{W}_{\otimes l}^\top)^{\frac{1}{2}} \right)^{\frac{1}{2}} \right) \right)^{\frac{1}{2}} \\ &= \sqrt{\|\mathbf{W}_{\otimes l}(\mathbf{W}_{\otimes L}^\top - \mu_0)\|^2 + \frac{(\sqrt{n-1} - \sqrt{n})^2 \sigma_0^2}{n} \|\mathbf{W}_{\otimes l}\|_{\text{F}}^2}. \end{aligned}$$

Similarly, we have

$$\mathcal{W}_2(P_{X_i|Y_i, \mathbf{W}, \mathbf{W}_{\otimes L}}, P_{X|Y, \mathbf{W}, \mathbf{W}_{\otimes L}}) = \sqrt{\|(\mathbf{W}_{\otimes L}^\top - \mu_0)\|^2 + \frac{d_0(\sqrt{n-1} - \sqrt{n})^2 \sigma_0^2}{n}}.$$

Since the activation functions are all 1-Lipschitz, i.e.,  $\rho_l = 1$  for all  $l = 1, \dots, L$  and the loss function  $\tilde{\ell}$  is  $4\sqrt{2}$ -Lipschitz, we have  $\bar{\rho}_l(\mathbf{W}) = 4\sqrt{2}(1 \vee \prod_{j=l+1}^L \|\mathbf{W}_j\|_{\text{op}})$  for  $l = 0, 1, \dots, L$ .

For notational simplicity, let  $\mathbf{W}_{\otimes 0} = \mathbf{W}_0 = \mathbf{I}_{d_0}$  and  $r_0 = d_0$ . Here we use  $\|\cdot\|_{\text{F}}$  of a vector to equivalently denote its Euclidean norm, with a slight abuse of notations. Then the generalization error is upper bounded by

$$\begin{aligned} &\text{gen}(P_{\mathbf{W}|D_n}, P_{X,Y}) \\ &\leq \min \left\{ \min_{l=1, \dots, L} \mathbb{E} \left[ \bar{\rho}_l(\mathbf{W}) \sqrt{\|\mathbf{W}_{\otimes l}(\mathbf{W}_{\otimes L}^\top - \mu_0)\|^2 + \frac{(\sqrt{n-1} - \sqrt{n})^2 \sigma_0^2}{n} \|\mathbf{W}_{\otimes l}\|_{\text{F}}^2} \right], \right. \\ &\quad \left. \mathbb{E} \left[ \bar{\rho}_0(\mathbf{W}) \sqrt{\|(\mathbf{W}_{\otimes L}^\top - \mu_0)\|^2 + \frac{d_0(\sqrt{n-1} - \sqrt{n})^2 \sigma_0^2}{n}} \right] \right\} \\ &\stackrel{(a)}{\leq} \min \left\{ \min_{l=1, \dots, L} \mathbb{E} \left[ \bar{\rho}_l(\mathbf{W}) \left( \|\mathbf{W}_{\otimes l}(\mathbf{W}_{\otimes L}^\top - \mu_0)\| + \frac{(\sqrt{n} - \sqrt{n-1})\sigma_0}{\sqrt{n}} \|\mathbf{W}_{\otimes l}\|_{\text{F}} \right) \right], \right. \\ &\quad \left. \mathbb{E} \left[ \bar{\rho}_0(\mathbf{W}) \left( \|(\mathbf{W}_{\otimes L}^\top - \mu_0)\| + \frac{\sqrt{d_0}(\sqrt{n} - \sqrt{n-1})\sigma_0}{\sqrt{n}} \right) \right] \right\} \\ &\stackrel{(b)}{\leq} \min \left\{ \min_{l=1, \dots, L} \mathbb{E} \left[ \bar{\rho}_l(\mathbf{W}) \|\mathbf{W}_{\otimes l}\|_{\text{F}} \left( \|(\mathbf{W}_{\otimes L}^\top - \mu_0)\| + \frac{(\sqrt{n} - \sqrt{n-1})\sigma_0}{\sqrt{n}} \right) \right], \right. \\ &\quad \left. \mathbb{E} \left[ \bar{\rho}_0(\mathbf{W}) \sqrt{d_0} \left( \|(\mathbf{W}_{\otimes L}^\top - \mu_0)\| + \frac{(\sqrt{n} - \sqrt{n-1})\sigma_0}{\sqrt{n}} \right) \right] \right\} \\ &= \min_{l=0, \dots, L} \mathbb{E} \left[ \bar{\rho}_l(\mathbf{W}) \|\mathbf{W}_{\otimes l}\|_{\text{F}} \left( \|\mathbf{W}_{\otimes L}^\top - \mu_0\| + \frac{(\sqrt{n} - \sqrt{n-1})\sigma_0}{\sqrt{n}} \right) \right] \\ &\stackrel{(c)}{\leq} \min_{l=0, \dots, L} \mathbb{E} [\bar{\rho}_l(\mathbf{W})^2 \|\mathbf{W}_{\otimes l}\|_{\text{F}}^2]^{\frac{1}{2}} \mathbb{E} \left[ \left( \|(\mathbf{W}_{\otimes L}^\top - \mu_0)\| + \frac{(\sqrt{n} - \sqrt{n-1})\sigma_0}{\sqrt{n}} \right)^2 \right]^{\frac{1}{2}} \\ &\stackrel{(d)}{\leq} \min_{l=0, \dots, L} \mathbb{E} [\bar{\rho}_l(\mathbf{W})^2 \|\mathbf{W}_{\otimes l}\|_{\text{F}}^2]^{\frac{1}{2}} \left( \frac{\sqrt{d_0}\sigma_0}{\sqrt{n}} + \frac{(\sqrt{n} - \sqrt{n-1})\sigma_0}{\sqrt{n}} \right) \\ &= \frac{4\sqrt{2}\sigma_0(\sqrt{d_0} + (\sqrt{n} - \sqrt{n-1}))}{\sqrt{n}} \min_{l=0, \dots, L} \mathbb{E} \left[ \left( 1 \vee \prod_{j=l+1}^L \|\mathbf{W}_j\| \right)^2 \|\mathbf{W}_{\otimes l}\|_{\text{F}}^2 \right]^{\frac{1}{2}}, \end{aligned}$$

where (a) follows since  $\sqrt{a^2 + b^2} \leq |a| + |b|$ , (b) follows from the Cauchy-Schwarz inequality, (c) follows from the Hölder's inequality, and (d) follows from Minkowski's inequality and  $\frac{n}{\sigma_0^2} \|\mathbf{W}_{\otimes L}^\top - \mu_0\|^2 \sim \chi_{d_0}^2$ .  $\square$

*Proof of Example 1.* Since  $\mathbf{W}_l$  is  $(2 \times 2)$  rotation matrix multiplied by a scalar factor  $C_l$  and  $\mathbf{W}_L = (0, C_L)$  is a row vector,  $\|\mathbf{W}_l\| = C_l$  for  $l = 1, \dots, L$ . We have  $\bar{\rho}_l(\mathbf{W}) = 4\sqrt{2}(1 \vee \prod_{j=l+1}^L \|\mathbf{W}_j\|) = 4\sqrt{2}(1 \vee \prod_{j=l+1}^L C_j)$  for  $l = 0, 1, \dots, L$ .  $\square$